
SAMPLE SIZE ESTIMATION FOR CHEST X-RAY CLASSIFICATION WITH FOUNDATION MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

The integration of deep learning models into clinical practice, particularly in radiology, is often hindered by the need for large, meticulously labeled datasets, which entails significant time and financial costs. While foundation models substantially reduce this dependency, a critical question remains: what is the minimum amount of annotated data sufficient to achieve clinically acceptable accuracy? In this work, we introduce a methodology for accurately predicting sample size requirements by modeling learning curves with a power law. Our study demonstrates that modern foundation models, such as XrayCLIP and XraySigLIP, not only outperform traditional architectures but also achieve high ROC-AUC scores with significantly fewer training examples. A key finding of our research is the evidence that the learning dynamics observed with a sample of just 50 labeled cases can predict the model’s asymptotic performance with high precision. Thus, our study offers a scientifically grounded approach to optimizing the data annotation process, enabling researchers and clinicians to minimize costs and accelerate the development of reliable diagnostic tools.

1 INTRODUCTION

While significant literature exists on deep learning methods for chest X-ray classification (Mee-deniyi et al., 2022), comparatively little attention has been paid to efficient training size estimation in this context (Viering & Loog, 2022). Recent progress in foundation models has further heightened the importance of this question: not only can these models achieve higher accuracy, but their learning curves may also be more predictable with fewer labeled samples.

A key principle of this study is the rigorous and unbiased evaluation of labeling efficiency. Although the MIMIC-CXR dataset contains structured labels for some pathologies, we confirm that none of the evaluated foundation models—RAD-DINO-MAIRA-2, XrayCLIP, and XraySigLIP—were pre-trained using these specific structured annotations. The RAD-DINO-MAIRA-2 model was trained solely on images, while XrayCLIP and XraySigLIP were trained on image-text pairs using unstructured, free-form reports. This methodological choice is crucial as it ensures a level playing field for all pathologies in terms of the models’ prior knowledge, allowing for a fair assessment of their sample efficiency in a real-world fine-tuning scenario.

Motivated by these considerations, we propose a systematic approach to estimate how many annotated examples a given model requires to meet a clinically relevant ROC-AUC threshold, leveraging power-law fitting of the learning curves.

2 RELATED WORK

Chest X-ray is a crucial diagnostic imaging modality that provides rapid and cost-effective insights into various pulmonary and cardiac conditions (Raouf et al., 2012). The classification of chest X-ray pathologies is well-studied, and training machine learning models for new conditions is relatively straightforward, although it still relies on large annotated datasets to achieve clinically acceptable accuracy (Çallı et al., 2021).

Recently, general self-supervised learning (SSL) frameworks such as DINO (Oquab et al., 2023) and CLIP (Radford et al., 2021) have shown great promise for imaging tasks by learning robust feature

054 representations from massive unlabeled datasets. These frameworks differ in their training objec-
055 tives: DINO is purely image-based self-supervision, whereas CLIP leverages paired text–image data
056 for multi-modal alignment. Building on these advances, specialized chest X-ray foundation models
057 (*e.g.*, RadDINO (Pérez-García et al., 2024), XraySigLIP/XrayCLIP (Chen et al., 2024)) adapt these
058 frameworks to chest x-ray imaging.

059 Beyond predicting performance scaling, another prominent approach to data-efficient learning is
060 to improve the quality of the learned representations themselves. For instance, Supervised Con-
061 trastive Learning (SupCon) has emerged as a powerful technique that goes beyond the standard
062 cross-entropy loss function (Moradinasab et al., 2024). SupCon aims to pull representations of ex-
063 amples from the same class closer together in the latent space while pushing apart examples from
064 different classes (Wu et al., 2023). This approach has shown significant promise in medical imag-
065 ing, where it can help learn more discriminative features even with limited labeled data. However,
066 SupCon can face challenges such as "class collapse", where intra-class variance is lost (Chen et al.,
067 2022), and its effectiveness can diminish under severe class imbalance, a common scenario in med-
068 ical diagnostics. While these methods aim to fundamentally alter the training objective to learn a
069 *better* representation space, our work addresses the complementary, pragmatic question: for widely-
070 used, off-the-shelf foundation models, how can one reliably predict their performance trajectory
071 to inform annotation budgets? Thus, our study focuses on the practical estimation of sample size
072 requirements rather than the development of new representation learning techniques.

073 Though no works have directly explored learning curve estimates for chest X-ray classification tasks,
074 many investigations in other domains (*e.g.*, machine translation, image recognition, and speech
075 recognition) rely on power-law approximations to characterize how performance improves as the
076 training set size grows (Cortes et al., 1993; Gu et al., 2001; Hestness et al., 2017). Nonetheless,
077 numerous studies reveal that learning curves can be well-behaved or ill-behaved, with phenomena
078 such as double descent and peaking complicating straightforward sample-size extrapolation (Raudys
079 & Duin, 1998; Devroye et al., 2013; Nakkiran et al., 2020). A popular strategy is to estimate the
080 asymptotic accuracy by measuring the early slope of a power-law fit and extrapolating the eventual
081 plateau in performance (Hoiem et al., 2021; Frey & Fisher, 1999; Kolachina et al., 2012). Addi-
082 tionally, a relevant idea is to incorporate progressive sampling, dynamically refining the power-law
083 estimate of the learning curve so as to reduce annotation overhead (Provost et al., 1999).

085 3 METHODS

087 3.1 DATASET CONSTRUCTION

089 A popular open dataset MIMIC-CXR (Johnson et al., 2019) was used as the data source for this
090 study. Using RadGraph annotations (Jain et al., 2021), we extracted structured “organ–pathology”
091 labels. These labels underwent a normalization to merge synonymous anatomical terms into uni-
092 fied categories (*e.g.*, unifying the tokens “lung” and “lobe”) and then split into two groups: normal
093 and pathological. Pathological annotations were further clustered according to their common patho-
094 genetic mechanisms to reduce redundancy. In total, 21 distinct pathology classes were selected.
095 An example of chest X-ray, corresponding RadGraph findings, and selected pathologies are shown
096 in Figure 1.

097 For each resulting pathology, we created a binary classification dataset, where a confirmed pathol-
098 ogy was labeled as a positive class. The negative class consisted of studies corresponding to a
099 normal anatomical-physiological state of the target organ, in a 1:5 ratio. If for some classes negative
100 examples were insufficient, existing data were duplicated to maintain the balance.

101 Each pathology-specific dataset was split into training, validation, and test subsets using a determi-
102 nistic method with a fixed seed value. The validation and hold-out test subsets were each assigned
103 10% of the total data in a stratified manner, preserving the 1:5 class ratio. The remaining 80% made
104 up the full training pool.

105 The choice of fixed increments for positive cases (ranging from 5 to 1000 samples) and a con-
106 stant 1:5 class ratio was a deliberate methodological decision. This controlled environment, while
107 an abstraction of clinical prevalence, was necessary to isolate the effect of the number of positive

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

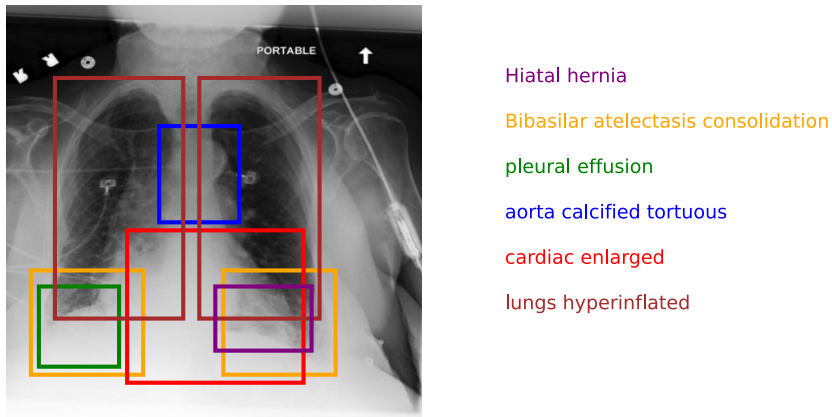


Figure 1: A chest X-ray with pathology labels extracted via RadGraph.

examples on performance and to ensure a systematic, comparable analysis across a wide range of pathologies and their varying data availability.

In each experiment below, the training sets were formed for a feature-based transfer learning regime as follows:

1. The number of positively annotated samples was purposely restricted with a value N_{cases} taken from the set $\{5, 10, 15, \dots, 45, 50, 100, 250, 500, 1000\}$.
2. From the initial training pool, N_{cases} pathological studies were randomly selected (using a unique seed for each experiment).
3. Negative cases were added in a 1:5 ratio to preserve the original balance.

The validation and the testing sets were the same for each pathology in all experiments.

3.2 MODELS TESTED

For feature extraction, we employed 3 different chest x-ray foundation models. The first, RAD-DINO-MAIRA-2 (Bannur et al., 2024), is a transformer pre-trained using the DINOv2 framework on a heterogeneous corpus of 1.2 million medical images. The second and the third are XraySigLIP and XrayCLIP – also transformer models, but pre-trained using the CLIP (Radford et al., 2021) and SigLIP (Zhai et al., 2023) frameworks, respectively, on a million image-text pairs from CheXinstruct (Chen et al., 2024) dataset. To verify the robustness of the results and to establish a baseline, we also used ResNet-50 convolutional neural network (He et al., 2015) pre-trained on ImageNet as a baseline encoder. For each of the tested models we constructed the classifier head by applying a dropout layer ($p = 0.1$) to the encoder’s pooled features, followed by a linear projection to a single output unit.

3.3 TRAINING PROCEDURE

Training was done using the `transformers` library (PyTorch backend) with the AdamW optimizer (binary cross-entropy loss, an initial learning rate 2×10^{-5}), a cosine annealing learning rate scheduler without warm-up, a batch size of 64, and an early stopping after 4 consecutive epochs without improvement in the validation loss. During training, the image encoder weights were frozen to retain their pre-trained representations, and only the linear binary classifier head was trained.

We applied train augmentations combining geometric and photometric modifications: a horizontal flip (50% probability), affine transformations with a random rotation between -90° and 90° and a rotation center is center of the image size, photometric adjustments via linear brightness adaptation

within $\pm 35\%$ of the original values alongside non-linear gamma correction of contrast in the same range, and spatial cropping with random square crops, covering 3–33% of the original image size.

3.4 POWER LAW FITTING

To model the scaling behavior of the classifier, we fit a power law function to the area under the receiver-operating characteristic curve (ROC-AUC) in the following form:

$$\text{ROC_AUC}(n) = \alpha - \frac{\beta}{n^\gamma}, \quad (1)$$

where n is the number of distinct positive examples in the training set, α represents the asymptotic performance, β controls the deviation from the asymptote, and γ governs the rate of convergence. Several curve-fitting approaches (linear, exponential, and power-law) were considered. The three-parameter power-law function, shown in Equation 1, was selected due to its consistently superior fit to the empirical learning curves across the range of pathologies and models evaluated, consistent with prior research on scaling laws in deep learning (Hestness et al., 2017; Kaplan et al., 2020; Viering & Loog, 2022).

To estimate the parameters α , β , and γ , we employ non-linear least squares fitting using the `curve_fit` function from `scipy.optimize`. For the fitting procedure, we specify an initial guess and bounds for the parameters to ensure reasonable behavior of the model. In our case, we set the initial guess as $\alpha = 0.95, \beta = 0.5, \gamma = 1.0$ with the following bounds: $\alpha \in [0.8, 1.0]$ ensuring the asymptotic value is near 1, $\beta \geq 0$ to maintain non-negative deviation, $\gamma \geq 0$ for a proper convergence rate.

For the fitted power law we use the following notation $\text{ROC_AUC}_{N_{\text{cases}}}(n)$ where N_{cases} is the maximum number of examples used to fit the curve. For example, $\text{ROC_AUC}_{20}(n)$ stands for the power law curve, fitted on the experimental data points $N_{\text{cases}} = 5, 10, 15, 20$. Finally, given the fitted curve, we draw a conclusion about the optimal number of required labeled samples n_o by evaluating where the curve starts exceeding a certain clinically-relevant threshold ($\text{ROC_AUC}_{N_{\text{cases}}}(n_o) = 90\%$)¹.

4 RESULTS

4.1 ALL PATHOLOGIES DATA POINTS

The results for all pathologies are presented in Table 1. For the 4 models and each of the pathologies we provide the experimental ROC-AUC on all the training data available for this pathology, and the expected number of cases needed to reach ROC-AUC 0.9. This number was calculated by fitting a power law to all the experimental data using less than 50 training samples and using it to calculate the number of examples needed.

4.2 NUMBER OF TRAINING EXAMPLES VS EXPERIMENTAL ROC-AUC

Figure 2 illustrates the core principle of our methodology using the “lobe mass” pathology as an example. It demonstrates that power-law curves fitted on a small fraction of the data (*e.g.*, ROC_AUC_{50} , fitted on $N \leq 50$ cases) closely track the true empirical learning curve ($\text{ROC_AUC} \pm 1 \text{ Std}$). This visual evidence supports our hypothesis that early learning dynamics are highly predictive of future performance, forming the basis for reliable sample size estimation.

4.3 COMPARISON OF EARLY SLOPE AND FINAL PERFORMANCE

To quantify the relationship between early-stage learning behavior and final performance, we calculated the Pearson correlation coefficient (r) between the initial slope of the learning curve and

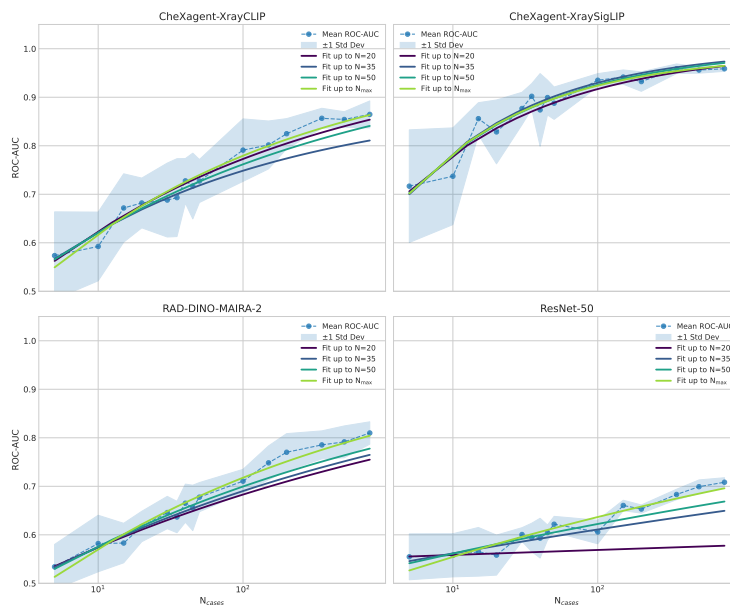
¹The ROC-AUC threshold of 0.90 used throughout this study was selected as an illustrative benchmark for simplicity and consistency. However, our methodology is generalizable and can readily accommodate any clinically relevant performance threshold, allowing practitioners to adjust the labeling requirements according to specific diagnostic standards.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241

Table 1: Performance metrics with best values highlighted in bold. Best ROC-AUC and best n@90 are shown in bold.

Pathology	ResNet-50		RAD-DINO		XrayCLIP		XraySigLIP	
	roc	n@90	roc	n@90	roc	n@90	roc	n@90
pulmonary_fibrosis	0.85	2545	0.92	104	0.97	24	0.99	8
pericardial_effusion	0.65	>1M	0.73	98922	0.77	5486	0.92	79
aortic_dissection	0.72	>1M	0.81	241	0.71	4656	0.98	18
hiatal_hernia	0.78	>1M	0.92	646	0.90	317	0.93	120
lobe_mass	0.71	>1M	0.81	156K	0.86	7605	0.96	53
hemidiaphragm_eventration	0.78	1805	0.84	5315	0.86	688	0.84	8954
fissure_fluid	0.64	inf	0.67	162K	0.75	246K	0.95	45
spine_deformities	0.81	1159	0.75	>1M	0.87	>1M	0.91	77
pulmonary_hypertension	0.56	>1M	0.72	>1M	0.70	>1M	0.86	1461
clavicular_fracture	0.56	>1M	0.69	>1M	0.74	172K	0.73	>1M
esophagus_dilated	0.45	inf	0.57	>1M	0.63	>1M	0.82	161
lung_edema	0.85	510	0.77	107K			1.00	15
diaphragms_flattened	0.85	412	0.85	13228	0.93	233	0.90	272
rib_fractures	0.60	>1M	0.73	>1M	0.89	999	0.87	92
lung_aeration	0.67	>1M	0.72	229K	0.96	49		
hilar_mass	0.69	>1M	0.80	114K	0.91	306	0.91	80
aorta_calcification	0.74	>1M	0.75	267K	0.88	338	0.89	97
mediastinum_shift	0.68	>1M	0.82	7904	0.95	18	0.98	22
lung_atelectasis	0.64	>1M	0.58	>1M	0.89	7071	0.81	67
pleural_air	0.75	>1M	0.85	>1M	0.96	22	0.95	34
cardiac_enlarged	0.63	>1M	0.71	>1M	0.86	4439	0.86	2365

242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261



262
263
264
265
266
267
268
269

Figure 2: ROC-AUC vs the number of training examples for Lobe mass pathology.

the final achieved ROC-AUC. The slope, defined as the derivative of the fitted power-law function ($ROC_AUC'(n) = \frac{\beta\gamma}{n^{\gamma+1}}$), was evaluated at a small sample size ($n = 5$). As shown in Figure 3, we observe a strong positive correlation that increases as more initial data points are used for fitting (e.g., for ResNet-50, r increases from 0.42 to 0.84 as the data cutoff for fitting is raised from 10 to 50 cases). This strong correlation provides statistical validation for our central finding: **the steepness of the initial learning curve is a reliable predictor of the model’s final performance plateau.**

270
271
272
273
274
275
276
277
278
279
280
281
282
283

This allows for confident extrapolation from small, low-cost pilot experiments to estimate the full annotation budget required to meet a target performance metric.

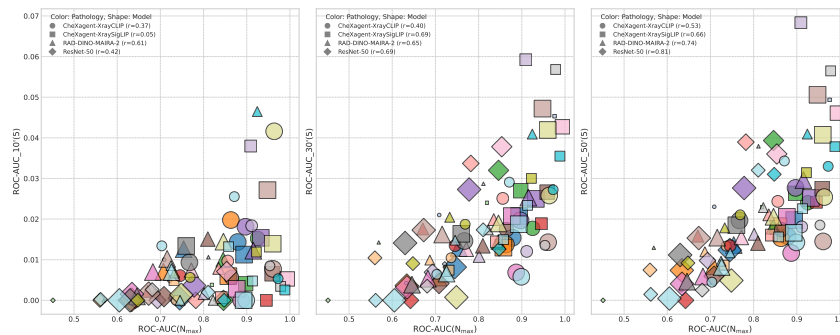


Figure 3: Correlation between the derivatives of the fitted ROC-AUC at $n=5$ and the value of ROC-AUC(N_{max}).

284
285
286
287

4.4 ERROR IN PREDICTED VS. OBSERVED PLATEAU

288
289
290
291
292
293
294

Beyond measuring correlation, we also assessed absolute prediction error in estimating the final ROC-AUC. For each model-pathology pair, we used the power-law curve fitted at $N_{cases} = 20$ and $N_{cases} = 40$ to extrapolate to N_{max} equal to the maximum number of examples for this pathology. We then compared this predicted $ROC_AUC_{N_{cases}}(N_{max})$ with the actual measured value $ROC_AUC(N_{max})$.

295
296
297

Figure 4 depicts how the mean absolute error (MAE) across pathologies and models evolves as we gradually increase the cutoff for fitting. Notably, the MAE decreases rapidly up to about 50-100 labeled cases, after which the benefit of additional data for partial fits diminishes.

298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313

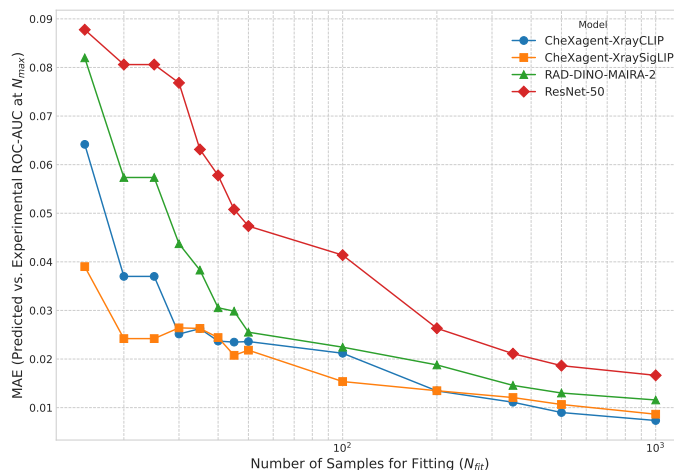


Figure 4: MAE between experimental ROC-AUC and ROC-AUC predicted on limited number of training examples.

314
315
316
317

5 DISCUSSION

318
319
320
321
322
323

In our experiments, we demonstrate that a straightforward process of running multiple training subsets (5 to 50 positive cases, with negative samples at a fixed ratio) allows for reliable fitting of power-law curves. By examining the initial slope and partial plateaus of these curves, we can extrapolate the performance of fine-tuned foundation models for higher training sizes. This approach is particularly useful in real-world settings where annotation is expensive, since it indicates when

324 additional labeling provides diminishing returns. Our findings show that, in many cases, labeling on
325 the order of 50 to 100 positive samples per pathology is sufficient to predict—and often achieve—
326 competitive diagnostic accuracy levels.

327 We acknowledge several deliberate limitations in the scope of this study. Our analysis was restricted
328 to binary classification for 21 pathologies from a single large dataset. This focus was a methodolog-
329 ical choice to ensure a clear, controlled, and reproducible investigation of scaling laws, avoiding the
330 complex confounding factors of multi-class classification or dataset shift. While the principles of
331 power-law fitting are general, future work is needed to validate the specific scaling coefficients and
332 predictive accuracy of this method across different datasets, imaging modalities, and in multi-class
333 scenarios.

334 Moreover, foundation models consistently outperform the conventional ResNet-50 baseline, under-
335 scoring not only their superior accuracy but also the improved predictability of their learning curves
336 from limited data. Their higher initial performance effectively reduces both the total labels needed
337 and the associated clinical costs.

338 The ROC-AUC threshold of 0.90 was used in this study as a consistent and illustrative benchmark.
339 However, our methodology is fundamentally target-agnostic. Practitioners can readily substitute
340 any clinically relevant performance threshold—be it sensitivity, specificity, or a different ROC-AUC
341 value—to tailor the sample size estimation to their specific diagnostic needs and regulatory stan-
342 dards. This adaptability is key to the practical utility of our framework. We anticipate that this
343 approach—train on subsets, fit a power law, then extrapolate to an ROC-AUC target—will inform
344 practitioners attempting to balance annotation budgets with diagnostic performance demands when
345 deploying chest X-ray classifiers for new pathologies.

347 REFERENCES

- 349 Shruthi Bannur, Kenza Bouzid, Daniel C Castro, Anton Schwaighofer, Anja Thieme, Sam Bond-
350 Taylor, Maximilian Ilse, Fernando Pérez-García, Valentina Salvatelli, Harshita Sharma, et al.
351 Maira-2: Grounded radiology report generation. *arXiv preprint arXiv:2406.04449*, 2024.
- 352 Erdi Çallı, Ecem Sogancioglu, Bram van Ginneken, Kicky G van Leeuwen, and Keelin Murphy.
353 Deep learning for chest x-ray analysis: A survey. *Medical image analysis*, 72:102125, 2021.
- 354 Mayee F Chen, Daniel Y Fu, Avanika Narayan, Michael Zhang, Zhao Song, Kayvon Fatahalian, and
355 Christopher Ré. Perfectly balanced: Improving transfer and robustness of supervised contrastive
356 learning. In *International Conference on Machine Learning*, pp. 3069–3091. PMLR, 2022.
- 357 Zhihong Chen, Maya Varma, Jean-Benoit Delbrouck, Magdalini Paschali, Louis Blankemeier,
358 Dave Van Veen, Jeya Maria Jose Valanarasu, Alaa Youssef, Joseph Paul Cohen, Eduardo Pontes
359 Reis, Emily B. Tsai, Andrew Johnston, Cameron Olsen, Tanishq Mathew Abraham, Sergios Gat-
360 tidis, Akshay S Chaudhari, and Curtis Langlotz. Chexagent: Towards a foundation model for
361 chest x-ray interpretation. *arXiv preprint arXiv:2401.12208*, 2024.
- 362 Corinna Cortes, Lawrence D Jackel, Sara Solla, Vladimir Vapnik, and John Denker. Learning
363 curves: Asymptotic values and rate of convergence. In *Advances in neural information processing*
364 *systems*, volume 6, 1993.
- 365 Luc Devroye, László Györfi, and Gábor Lugosi. *A probabilistic theory of pattern recognition*, vol-
366 ume 31. Springer Science & Business Media, 2013.
- 367 Lewis J Frey and Douglas H Fisher. Modeling decision tree performance with the power law. In
368 *Seventh international workshop on artificial intelligence and statistics*. PMLR, 1999.
- 369 Baohua Gu, Feifang Hu, and Huan Liu. Modelling classification performance for large data sets:
370 An empirical study. In *Advances in Web-Age Information Management: Second International*
371 *Conference, WAIM 2001 Xi'an, China, July 9–11, 2001 Proceedings 2*, pp. 317–328. Springer,
372 2001.
- 373 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
374 nition. *arXiv preprint arXiv:1512.03385*, 2015.

378 Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad,
379 Md Mostofa Ali Patwary, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable,
380 empirically. *arXiv preprint arXiv:1712.00409*, 2017.

381

382 Derek Hoiem, Tanmay Gupta, Zhizhong Li, and Michal Shlapentokh-Rothman. Learning curves
383 for analysis of deep networks. In *International conference on machine learning*, pp. 4287–4296.
384 PMLR, 2021.

385 Saahil Jain, Ashwin Agrawal, Adriel Saporta, Steven QH Truong, Du Nguyen Duong, Tan Bui,
386 Pierre Chambon, Yuhao Zhang, Matthew P Lungren, Andrew Y Ng, et al. Radgraph: Extracting
387 clinical entities and relations from radiology reports. *arXiv preprint arXiv:2106.14463*, 2021.

388

389 Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz, Nathaniel R Greenbaum, Matthew P Lun-
390 gren, Chih-ying Deng, Roger G Mark, and Steven Horng. Mimic-cxr, a de-identified publicly
391 available database of chest radiographs with free-text reports. *Scientific data*, 6(1):317, 2019.

392 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,
393 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
394 models. *arXiv preprint arXiv:2001.08361*, 2020.

395

396 Prasanth Kolachina, Nicola Cancedda, Marc Dymetman, and Sriram Venkatapathy. Prediction of
397 learning curves in machine translation. In *Proceedings of the 50th Annual Meeting of the Associ-
398 ation for Computational Linguistics (Volume 1: Long Papers)*, pp. 22–30, 2012.

399 Dulani Meedeniya, Hashara Kumarasinghe, Shammi Kolonne, Chamodi Fernando, Isabel De la
400 Torre Díez, and Goncalo Marques. Chest x-ray analysis empowered with deep learning: A sys-
401 tematic review. *Applied Soft Computing*, 126:109319, 2022.

402

403 Nazanin Moradinasab, Suchetha Sharma, Ronen Bar-Yoseph, Shlomit Radom-Aizik, Kenneth C.
404 Bilchick, Dan M. Cooper, Arthur Weltman, and Donald E. Brown. Universal representa-
405 tion learning for multivariate time series using the instance-level and cluster-level supervised
406 contrastive learning. *Data Mining and Knowledge Discovery*, 38(3):1493–1519, 2024. doi:
407 10.1007/s10618-024-01006-1.

408

409 Preetum Nakkiran, Prayaag Venkat, Sham Kakade, and Tengyu Ma. Optimal regularization can
410 mitigate double descent. *arXiv preprint arXiv:2003.01897*, 2020.

411

412 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,
413 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning
414 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.

415

416 Fernando Pérez-García, Harshita Sharma, Sam Bond-Taylor, Kenza Bouzid, Valentina Salvatelli,
417 Maximilian Ilse, Shruthi Bannur, Daniel C Castro, Anton Schwaighofer, Matthew P Lungren,
418 et al. Rad-dino: Exploring scalable medical image encoders beyond text supervision. *arXiv
419 preprint arXiv:2401.10815*, 2024.

420

421 Foster Provost, David Jensen, and Tim Oates. Efficient progressive sampling. In *Proceedings of
422 the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp.
423 23–32, 1999.

424

425 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
426 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
427 models from natural language supervision. In *International conference on machine learning*, pp.
428 8748–8763. PmLR, 2021.

429

430 Suhail Raoof, David Feigin, Arthur Sung, Sabiha Raoof, Lavanya Irugulpati, and Edward C
431 Rosenow III. Interpretation of plain chest roentgenogram. *Chest*, 141(2):545–558, 2012.

432

433 Sarunas Raudys and Robert PW Duin. Expected classification error of the fisher linear classifier
434 with pseudo-inverse covariance matrix. *Pattern recognition letters*, 19(5-6):385–392, 1998.

435

436 Tom Viering and Marco Loog. The shape of learning curves: a review. *IEEE Transactions on
437 Pattern Analysis and Machine Intelligence*, 2022.

432 Yilei Wu, Zijian Dong, Chongyao Chen, Wangchunshu Zhou, and Juan Helen Zhou. Supremix:
433 Supervised contrastive learning for medical imaging regression with mixup. *arXiv preprint*
434 *arXiv:2309.16633*, 2023.

435 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language
436 image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*,
437 pp. 11975–11986, 2023.

438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485